

# Density Sensitive Hashing

Yue Lin, Deng Cai, *Member, IEEE*, Cheng Li

**Abstract**—Nearest neighbors search is a fundamental problem in various research fields like machine learning, data mining and pattern recognition. Recently, hashing-based approaches, e.g., Locality Sensitive Hashing (LSH), are proved to be effective for scalable high dimensional nearest neighbors search. Many hashing algorithms found their theoretic root in random projection. Since these algorithms generate the hash tables (projections) randomly, a large number of hash tables (i.e., long codewords) are required in order to achieve both high precision and recall. To address this limitation, we propose a novel hashing algorithm called *Density Sensitive Hashing* (DSH) in this paper. DSH can be regarded as an extension of LSH. By exploring the geometric structure of the data, DSH avoids the purely random projections selection and uses those projective functions which best agree with the distribution of the data. Extensive experimental results on real-world data sets have shown that the proposed method achieves better performance compared to the state-of-the-art hashing approaches.

**Index Terms**—Locality Sensitive Hashing, Random Projection, Clustering.

## 1 INTRODUCTION

Nearest Neighbors (NN) search is a fundamental problem and has found applications in many data mining tasks [9], [11], [14]. A number of efficient algorithms, based on pre-built index structures (e.g. KD-tree [4] and R-tree [2]), have been proposed for nearest neighbors search. Unfortunately, these approaches perform worse than a linear scan when the dimensionality of the space is high [5].

Given the intrinsic difficulty of exact nearest neighbors search, many hashing algorithms are proposed for Approximate Nearest Neighbors (ANN) search [1], [8], [10]. The key idea of these approaches is to generate binary codewords for high dimensional data points that preserve the similarity between them. Roughly, these hashing methods can be divided into two groups, the random projection based methods and the learning based methods.

Many hashing algorithms are based on the random projection, which has been proved to be an effective method to preserve pairwise distances for data points. One of the most popular methods is Locality Sensitive Hashing (LSH) [1], [8], [10], [12]. Given a database with  $n$  samples, LSH makes no prior assumption about the data distribution and offers probabilistic guarantees of retrieving items within  $(1 + \epsilon)$  times the optimal similarity, with query times that are sub-linear with respect to  $n$  [22], [27]. However, according to the Johnson Lindenstrauss Theorem [17], LSH needs  $O(\ln n / \epsilon^2)$  random projections to preserve the pairwise distances, where  $\epsilon$  is the relative error. Therefore, in order to increase the probability that similar objects are mapped to similar

hash codes, LSH needs to use many random vectors to generate the hash tables (a long codeword), leading to a large storage space and a high computational cost.

Aiming at making full use of the structure of the data, many learning-based hashing algorithms [6], [15], [16], [31], [32], [35], [38] are proposed. Most of these algorithms exploit the spectral properties of the data affinity (i.e., item-item similarity) matrix for binary coding. Despite the success of these approaches for relatively small codes, they often fail to make significant improvement as the code length increases [19].

In this paper, we propose a novel hashing algorithm called *Density Sensitive Hashing* (DSH) for effective high dimensional nearest neighbors search. Our algorithm can be regarded as an extension of LSH. Different from all the existing random projection based hashing methods, DSH tries to utilize the geometric structure of the data to guide the projections (hash tables) selection. Specifically, DSH uses  $k$ -means to roughly partition the data set into  $k$  groups. Then for each pair of adjacent groups, DSH generates one projection vector which can well split the two corresponding groups. From all the generated projections, DSH select the final ones according to the maximum entropy principle, in order to maximize the information provided by each bit. Experimental results show the superior performance of the proposed Density Sensitive Hashing algorithm over the existing state-of-the-art approaches.

The remainder of this paper is organized as follows. We introduce the background and review the related work in Section 2. Our Density Sensitive Hashing algorithm is presented in Section 3. Section 4 gives the experimental results that compared our algorithm with the state-of-the-art hashing methods on three real world large scale data sets. Conclusion remarks are provided in Section 5.

Y. Lin, D. Cai and C. Li are with the State Key Lab of CAD&CG, College of Computer Science, Zhejiang University, Hangzhou, Zhejiang, China, 310058. Email: linyue29@gmail.com, dengcai@cad.zju.edu.cn, licheng@zju.edu.cn.

## 2 BACKGROUND AND RELATED WORK

The generic hashing problem is the following. Given  $n$  data points  $\mathbf{X} = [\mathbf{x}_1, \dots, \mathbf{x}_n] \in \mathbb{R}^{d \times n}$ , find  $L$  hash functions to map a data point  $\mathbf{x}$  to a  $L$ -bits hash code

$$H(\mathbf{x}) = [h_1(\mathbf{x}), h_2(\mathbf{x}), \dots, h_L(\mathbf{x})],$$

where  $h_l(\mathbf{x}) \in \{0, 1\}$  is the  $l$ -th hash function. For the linear projection-based hashing, we have [33]

$$h_l(\mathbf{x}) = \text{sgn}(F(\mathbf{w}_l^T \mathbf{x} + t_l)) \quad (1)$$

where  $\mathbf{w}_l$  is the projection vector and  $t_l$  is the intercept. Different hashing algorithms aim at finding different  $F$ ,  $\mathbf{w}_l$  and  $t_l$  with respect to different objective functions.

One of the most popular hashing algorithms is Locality Sensitive Hashing (LSH) [1], [8], [10], [12]. LSH is fundamentally based on the random projection and uses randomly generated  $\mathbf{w}_l$ . The  $F$  in LSH is an identity function and  $t_l = 0$  for mean thresholding<sup>1</sup>. Thus, for LSH, we have

$$h_l(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}_l^T \mathbf{x} \geq 0 \\ 0 & \text{otherwise} \end{cases} \quad (2)$$

where  $\mathbf{w}_l$  is a vector generated from a zero-mean multivariate Gaussian  $\mathcal{N}(0, \mathbf{I})$  of the same dimension as the input  $\mathbf{x}$ . From the geometric point of view, the  $\mathbf{w}_l$  defines a hyperplane. The points on different sides of the hyperplane have the opposite labels. Using this hash function, two points' hash bits match with the probability proportional to the cosine of the angle between them [8]. Specifically, for any two points  $\mathbf{x}_i, \mathbf{x}_j \in \mathbb{R}^d$ , we have [22]:

$$\Pr[h_l(\mathbf{x}_i) = h_l(\mathbf{x}_j)] = 1 - \frac{1}{\pi} \cos^{-1}\left(\frac{\mathbf{x}_i^T \mathbf{x}_j}{\|\mathbf{x}_i\| \|\mathbf{x}_j\|}\right) \quad (3)$$

Based on this nice property, LSH have the probabilistic guarantees of retrieving items within  $(1 + \epsilon)$  times the optimal similarity, with query times that are sub-linear with respect to  $n$  [22], [27].

Empirical studies [1] showed that the LSH is significantly more efficient than the methods based on hierarchical tree decomposition. It has been successfully used in various applications in data mining [9], [14], computer vision [32], [34] and database [20], [21]. There are many extensions for LSH [18], [22], [25], [28]. Entropy based LSH [28] and Multi-Probe LSH [25], [18] are proposed to reduce the space requirement in LSH but need much longer time to deal with the query. The original LSH methods cannot apply for high-dimensional kernelized data when the underlying feature embedding for the kernel is unknown. To address this limitation, Kernelized Locality Sensitive Hashing is introduced in [22]. It suggests to approximate a normal distribution in the kernel space using only kernel comparisons [19]. In addition, the Shift Invariant Kernels Hashing [30], which is a distribution-free method based on the random features

mapping for shift-invariant kernels, is also proposed recently. This method has theoretical convergence guarantees and performs well for relatively large code sizes [13]. All these methods are fundamentally based on the random projection. According to the Jolson Lindenstrauss Theorem [17],  $O(\ln n / \epsilon^2)$  projective vectors are needed to preserve the pairwise distances of a database with size  $n$  for the random projection, where  $\epsilon$  is the relative error. Therefore, in order to increase the probability that similar objects are mapped to similar hash codes, these random projection based hashing methods need to use many random vectors to generate the hash tables (a long codeword), leading to a large storage space and a high computational cost.

To address the above limitation, many learning-based hashing methods [3], [6], [13], [15], [16], [19], [23], [24], [26], [27], [29], [31], [32], [35], [36], [37], [38] are proposed. PCA Hashing [34] might be the simplest one. It chooses  $\mathbf{w}_l$  in Eq.(1) to be the principal directions of data. Many other algorithms [24], [33], [35], [38] exploit the spectral properties of the data affinity (*i.e.*, item-item similarity) matrix for binary coding. The spectral analysis of the data affinity matrix is usually time consuming [7]. To avoid the high computational cost, Weiss *et al.* [35] made a strong assumption that data is uniformly distributed and proposed a Spectral Hashing method (SpH). The assumption in SpH leads to a simple analytical eigenfunction solution of 1-D Laplacians, but the geometric structure of the original data is almost ignored, leading to a suboptimal performance. Anchor Graph Hashing (AGH) [24] is a recently proposed method to overcome this shortcoming. AGH generates  $k$  anchor points from the data and represents all the data points by sparse linear combinations of the anchors. In this way, the spectral analysis of the data affinity can be efficiently performed. Some other learning based hashing methods include Semantic Hashing [31] which uses the stacked Restricted Boltzmann Machine (RBM) to generate the compact binary codes; Semi-supervised Sequential Projection Hashing (S3PH) [33] which can incorporate supervision information. Despite the success of these learning based hashing approaches for relatively small codes, they often fail to make significant improvement as the code length increases [19].

## 3 DENSITY SENSITIVE HASHING

In this section, we give the detailed description on our proposed *Density Sensitive Hashing* (DSH) which aims at overcoming the disadvantages of both random projection based and learning based hashing approaches. To guarantee the performance will increase as the code length increases, DSH adopts the similar framework as LSH. Different from LSH which generates the projections randomly, DSH uses the geometric structure of the data to guide the selection of the projections.

Figure 1 presents a toy example to illustrate the basic idea of our approach. There are four Gaussians in a

1. Without loss of generality, we assume that all the data points are centralized to have zero mean.

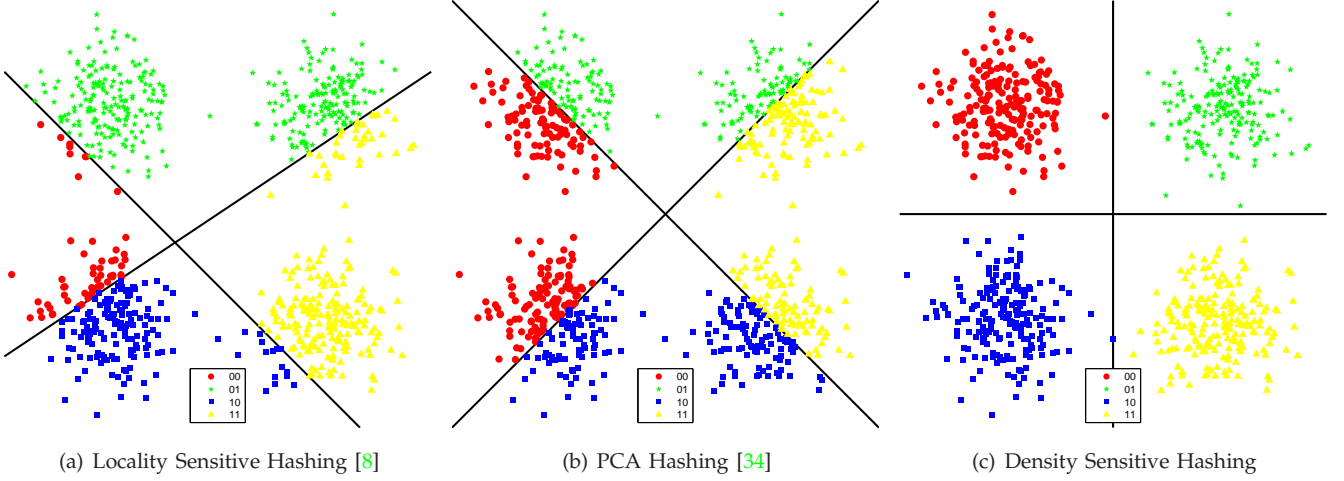


Fig. 1. Illustration of different hashing methods on a toy data set. There are four Gaussians in a two dimensional plane and one is asked to encode the data using two-bits hash codes. (a) LSH generates the projections randomly and it is very possible that the data points from the same Gaussian will be encoded by different hash codes. (b) PCA Hashing uses the principle directions of the data as the projective vectors. In our example, all the four Gaussians are split and PCA Hashing generates an unsatisfactory coding. (c) Considering the geometric structure of the data (density of the data), our DSH generates perfect binary codes for this toy example.

two dimensional plane and one is asked to encode the data using two-bits hash codes. LSH [8] generates the projections randomly and it is very possible that the data points from the same Gaussian will be encoded by different hash codes. PCA Hashing [34] uses the principle directions of the data as the projective vectors. In our example, all the four Gaussians are split and PCA Hashing generates an unsatisfactory coding. Considering the geometric structure of the data (density of the data), our DSH generates perfect binary codes for this toy example. The detailed procedure of DSH will be provided in the following subsections.

### 3.1 Minimum Distortion Quantization

The first step of DSH is quantization of the data. Recently, Paulevé *et al.* [29] show that a quantized preprocess for the data points can significantly improve the performance of the nearest neighbors search. Motivated by this result, we use the  $k$ -means algorithm, one of the most popular quantization approaches, to partition the  $n$  points into  $k$  ( $k < n$ ) groups.

Let  $\mathcal{S} = \{\mathcal{S}_1, \dots, \mathcal{S}_k\}$  denote a given quantization result. The *distortion*, also known as the *Sum of Squared Error* (SSE), can be used to measure the quality of the given quantization:

$$SSE = \sum_{p=1}^k \sum_{\mathbf{x} \in \mathcal{S}_p} \|\mathbf{x} - \mu_p\|^2 \quad (4)$$

The  $\mu_p$  is the representative point of the  $p$ -th group  $\mathcal{S}_p$ .

By noticing

$$\begin{aligned} \frac{\partial SSE}{\partial \mu_i} &= \frac{\partial}{\partial \mu_i} \sum_{p=1}^k \sum_{\mathbf{x} \in \mathcal{S}_p} \|\mathbf{x} - \mu_p\|^2 \\ &= \sum_{p=1}^k \sum_{\mathbf{x} \in \mathcal{S}_p} \frac{\partial}{\partial \mu_i} \|\mathbf{x} - \mu_p\|^2 \\ &= \sum_{\mathbf{x} \in \mathcal{S}_i} 2(\mathbf{x} - \mu_i) = 0, \quad i = 1, \dots, k, \end{aligned}$$

we have:

$$\mu_i = \frac{1}{|\mathcal{S}_i|} \sum_{\mathbf{x} \in \mathcal{S}_i} \mathbf{x}, \quad i = 1, 2, \dots, k \quad (5)$$

It indicates that in order to minimize the distortion, we can choose the center point as the representative point for each group.

There are two points that needed to be highlighted for the  $k$ -means quantization in our approach:

- 1) In large scale applications, it can be time consuming to wait the  $k$ -means converges. Naturally, we can stop the  $k$ -means after  $p$  iterations, where  $p$  is a parameter. We found that a small number of  $p$  is usually enough (usually 5). This will be discussed in the our experiments.
- 2) In real applications, we do not know which is the best group number  $k$ . It seems that the bigger the  $k$ , the better performance we will get. It is simply because the quantization will have smaller error with a large number of groups. However, a large number of groups could lead to high computational cost in the quantization step. As will be described in the next subsection, the number of

groups decides the maximum code length DSH can generate. Thus, we set

$$k = \alpha L \quad (6)$$

where  $L$  is the code length and  $\alpha$  is a parameter.

### 3.2 Density Sensitive Projections Generation

Now we have the quantization result denoted by  $k$  groups  $S_1, \dots, S_k$  and the  $i$ -th group has the center  $\mu_i$ . Instead of generating projections randomly as LSH does, our DSH tries to use this quantization result to guide the projections generating process.

We define the  $r$ -nearest neighbors matrix  $\mathbf{W}$  of the groups as follows:

**Definition 1:  $r$ -Nearest Neighbors Matrix  $\mathbf{W}$**  of the groups.

$$W_{ij} = \begin{cases} 1, & \text{if } \mu_i \in N_r(\mu_j) \text{ or } \mu_j \in N_r(\mu_i) \\ 0, & \text{otherwise} \end{cases} \quad (7)$$

where  $N_r(\mu_i)$  denotes the set of  $r$  nearest neighbors of  $\mu_i$ .

With this definition, we can then define  $r$ -adjacent groups:

**Definition 2:  $r$ -Adjacent Groups:** Group  $S_i$  and group  $S_j$  are called  $r$ -adjacent groups if and only if  $W_{ij} = 1$ . Instead of picking a random projection, it is more natural to pick those projections which can well separate two adjacent groups.

For each pair of adjacent groups  $S_i$  and  $S_j$ , DSH uses the median plane between the centers of adjacent groups as the hyperplane to separate points. The median plane is defined as follows:

$$(\mathbf{x} - \frac{\mu_i + \mu_j}{2})^T (\mu_i - \mu_j) = 0 \quad (8)$$

One can easily verify that the hash function associated with this plane is defined as follows:

$$h(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} \geq t \\ 0 & \text{otherwise} \end{cases} \quad (9)$$

where

$$\mathbf{w} = \mu_1 - \mu_2, \quad t = (\frac{\mu_1 + \mu_2}{2})^T (\mu_1 - \mu_2) \quad (10)$$

### 3.3 Entropy Based Projections Selection

Given  $k$  groups, the previous step can generate around  $\frac{1}{2}kr$  projections. Since  $k = \alpha L$ , our DSH generates  $\frac{1}{2}\alpha r L$  projections so far. Each projection will lead to one bit in the code and the usual setting of the parameters  $\alpha, r$  will make  $\frac{1}{2}\alpha r L > L$ . Thus, our DSH needs a projections selection step which aims at selecting  $L$  projections from the candidate set containing  $\frac{1}{2}\alpha r L$  projections.

From the information theoretic point of view, a "good" binary codes should maximize the information/entropy provided by each bit [38]. Using maximum entropy principle, a binary bit that gives balanced partitioning of the data points provides maximum information [32]. Thus,

we compute the entropy of each candidate projection and select the projections which can split the data most equally.

Assume we have  $m$  candidate projections  $\mathbf{w}_1, \mathbf{w}_2, \dots, \mathbf{w}_m$ . For each projection, the data points are separated into two sets and labeled with opposite bit. We denote these two partitions as  $\mathcal{T}_{i0}$  and  $\mathcal{T}_{i1}$ , respectively. The entropy  $\delta_i$  with respect to the projection  $\mathbf{w}_i$  can be computed as:

$$\delta_i = -P_{i0} \log P_{i0} - P_{i1} \log P_{i1} \quad (11)$$

where:

$$P_{i0} = \frac{|\mathcal{T}_{i0}|}{|\mathcal{T}_{i0}| + |\mathcal{T}_{i1}|}, \quad P_{i1} = \frac{|\mathcal{T}_{i1}|}{|\mathcal{T}_{i0}| + |\mathcal{T}_{i1}|} \quad (12)$$

In practice, the database can be very large and computing the entropy of each projection with respect to the entire database is time consuming. Thus, we estimate the entropy simply by using the group centers. For group center  $\mu_i$ , we assign a weight  $\nu_i$  based on the size of the group.

$$\nu_i = \frac{|\mathcal{S}_i|}{\sum_{p=1}^k |\mathcal{S}_p|} \quad (13)$$

We denote the two sets of group centers as  $\mathcal{C}_{i0}$  and  $\mathcal{C}_{i1}$ . Then  $P_{i0}$  and  $P_{i1}$  can be computed as:

$$P_{i0} = \sum_{s \in \mathcal{C}_{i0}} \nu_s, \quad P_{i1} = \sum_{t \in \mathcal{C}_{i1}} \nu_t \quad (14)$$

This simplification significantly reduces the time cost on the entropy calculation.

After obtaining the entry  $\delta_i$  for each  $\mathbf{w}_i$ , we sort them in descending order and use the top  $L$  projections for creating the  $L$ -bit binary codes, according to Eq.(9). The overall procedure of our DSH algorithm is summarized in Alg. 1.

### 3.4 Computational Complexity Analysis

Given  $n$  data points with the dimensionality  $d$ , the computational complexity of DSH in the training stage is as follows:

- 1)  $O(\alpha L p n d)$ :  $k$ -means with  $p$  iterations to generate  $\alpha L$  groups (Step 1 in Alg. 1).
- 2)  $O((\alpha L)^2(d+r))$ : Find all the  $r$ -adjacent groups (Step 2 in Alg. 1).
- 3)  $O(\alpha L r d)$ : For each pair of adjacent groups, generate the projection and the intercept (Step 3 in Alg. 1).
- 4) Compute the entropy for all the candidate projections needs  $O((\alpha L)^2 d r)$  (Step 4 in Alg. 1).
- 5) The top  $L$  projections can be found within  $O(\alpha L r \log(\alpha L r))$ . The binary codes for data points can be obtained in  $O(L n d)$  (Step 5 in Alg. 1).

Considering  $\alpha L r \ll n$ , the overall computational complexity of DSH training is dominated by the  $k$ -means clustering step which is  $O(\alpha L p n d)$ . It is clear that DSH scales linearly with respect to the number of samples in the database.



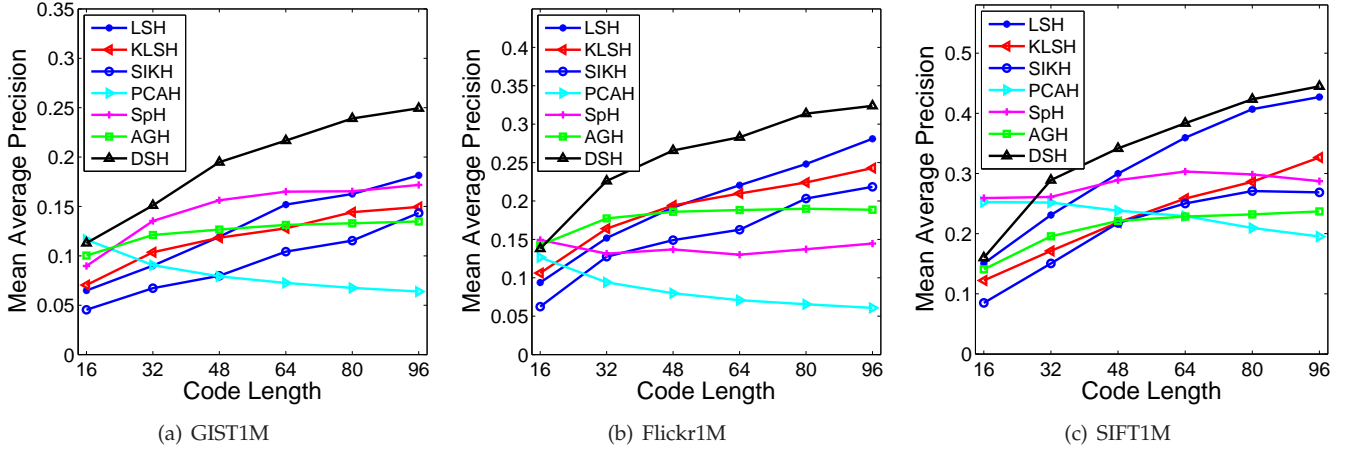


Fig. 2. The Mean Average Precision of all the algorithms on the three data sets.

---

#### Algorithm 1 Density Sensitive Hashing

---

##### Input:

- $n$  training samples  $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_n \in \mathbb{R}^d$ ;
- $L$ : the number of bits for hashing codes;
- $\alpha$ : the parameter controlling the groups number;
- $p$ : the number of iterations in the  $k$ -means;
- $r$ : the parameter for  $r$ -adjacent groups
- 1: Use  $k$ -means with  $p$  iterations to generate  $\alpha L$  groups, with centers  $\mu_1, \dots, \mu_{\alpha L}$ .
- 2: Generate the list of all  $r$ -adjacent groups based on the definition (1) and (2).
- 3: For each pair of adjacent groups, use Eq.(10) to generate the projection  $\mathbf{w}$  and intercept  $t$ .
- 4: Calculate the entropy of all the candidate projections using the weighted center points based on Eq.(11) and Eq.(14)
- 5: Sort the entropy values in descending order and use the top  $L$  projections to create binary codes according to Eq.(9).

##### Output:

The model:  $\{\mathbf{w}_i, t_i\}_{i=1}^L$   
 Binary hashing codes for the training samples:  $Y \in \{0, 1\}^{n \times L}$

---

In the testing stage, given a query point, DSH needs  $O(Ld)$  to compress the query point into a binary code, which is the same as the complexity of Locality Sensitive Hashing.

## 4 EXPERIMENT

In this section, we evaluate our DSH algorithm on the high dimensional nearest neighbor search problem. Three large scale real-world data sets are used in our experiments.

- **GIST1M**: It contains one million GIST features and each feature is represented by a 960-dim vector. This data set is publicly available<sup>2</sup>.

- **Flickr1M**: We collect one million images from the Flickr and use a feature extraction code<sup>3</sup> to extract a GIST feature for each image. Each image is represented by a 512-dim GIST feature vector. This data set is publicly available<sup>4</sup>.
- **SIFT1M**: It contains one million SIFT features and each feature is represented by a 128-dim vector. This data set is publicly available<sup>5</sup>.

For each data set, we randomly select 1k data points as the queries and use the remaining to form the gallery database. We use the same criterion as in [33], [36], that a returned point is considered to be a true neighbor if it lies in the top 2 percentile points closest (measured by the Euclidian distance in the original space) to the query. For each query, all the data points in the database are ranked according to their Hamming distances to the query. We evaluate the retrieval results by the Mean Average Precision (MAP) and the precision-recall curve [33]. In addition, we also report the training time and the testing time (the average time used for each query) for all the methods.

### 4.1 Compared Algorithms

Seven state-of-the-art hashing algorithms for high dimensional nearest neighbors search are compared as follows:

- Locality Sensitive Hashing (LSH) [8], which is based on the random projection. The projective vectors are randomly sampling from a  $p$ -stable distribution (e.g., Gaussian). We implement the algorithm by ourselves and make it publicly available<sup>6</sup>.
- Kernelized Locality Sensitive Hashing (KLSH) [22], which generalizes the LSH method to the kernel space. We use the code provided by the authors<sup>7</sup>.

3. <http://www.vision.ee.ethz.ch/~zhuji/felib.html>

4. <http://www.cad.zju.edu.cn/home/dengcai/Data/NNSDData.html>

5. <http://corpus-texmex.irisa.fr>

6. <http://www.cad.zju.edu.cn/home/dengcai/Data/DSH.html>

7. <http://www.cse.ohio-state.edu/~kulis/klsh/klsh.htm>

2. <http://corpus-texmex.irisa.fr>

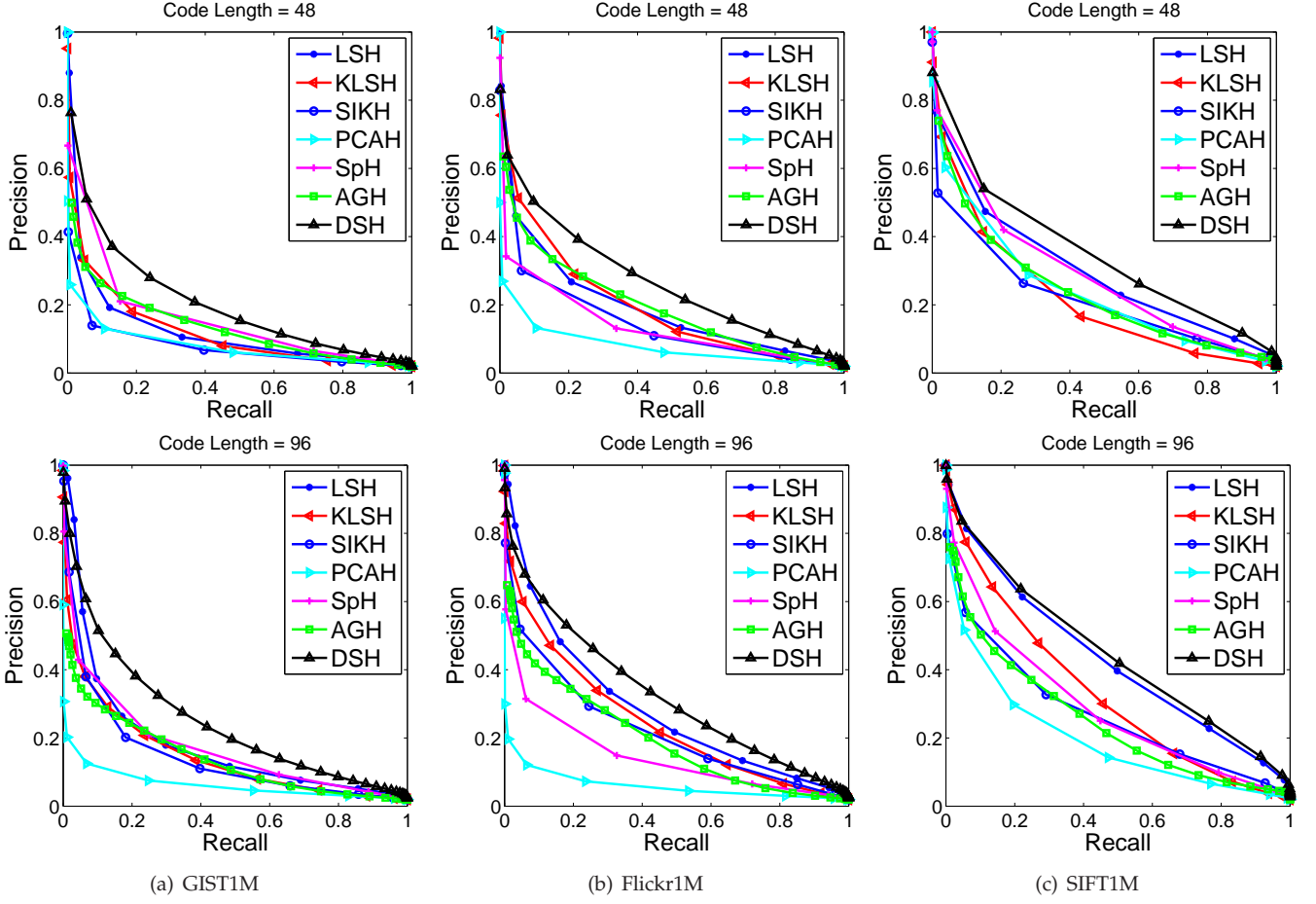


Fig. 3. The precision-recall curves of all algorithms on three data sets for the codes of 48 bits and 96 bits.

- Shift-Invariant Kernel Hashing (SIKH) [30], which is a distribution-free method based on the random features mapping for approximating shift-invariant kernels. The code is also publicly available<sup>8</sup>.
- Principle Component Analysis Hashing (PCAH) [34], which directly uses the top principal directions as the projective vectors to obtain the binary codes. The implementation of PCA is publicly available<sup>9</sup>.
- Spectral Hashing (SpH) [35], which is based on quantizing the values of analytical eigenfunctions computed along PCA directions of the data. We use the code provided by the authors<sup>10</sup>.
- Anchor Graph Hashing (AGH) [24], which constructs an anchor graph to speed up the spectral analysis procedure. AGH with two-layer is used in our comparison for its superior performance over AGH with one-layer [24]. We use the code provided by the authors<sup>11</sup> and the number of anchors is set to be 300 and the number of nearest neighbors is set to be 2 as suggested in [24].
- Density Sensitive Hashing (DSH), which is the

method introduced in this paper. For the purpose of reproducibility, we also make the code publicly available<sup>12</sup>. There are three parameters. We empirically set  $p = 3$  (the number of iterations in  $k$ -means),  $\alpha = 1.5$  (controlling the groups number),  $r = 3$  (for  $r$ -adjacent groups). A detailed analysis on the parameter selection will be provided later.

It is important to note that LSH, KLSH and SIKH are random projection based methods, while PCAH, SpH and AGH are learning based methods. Our DSH can be regarded as a combination of these two directions.

## 4.2 Experimental Results

Figure 2 shows the MAP curves of all the algorithms on the GIST1M, Flickr1M and SIFT1M data sets, respectively. We can see that the three random projection based methods (LSH, KLSH and SIKH) have a low MAP when the code length is short. As the code length increases, the performances of all the three methods consistently increases. On the other hand, the learning based methods (PCAH, SpH and AGH) have a high MAP when the code length is short. However, they fail to make significant improvements as the code length increases. Particularly, the performance of PCAH decreases as the code length

8. <http://www.unc.edu/~yunchao/itq.htm>

9. <http://www.cad.zju.edu.cn/home/dengcai/Data/DimensionReduction.htm>

10. <http://www.cs.huji.ac.il/~yweiss/SpectralHashing/>

11. <http://www.ee.columbia.edu/~wliu/>

12. <http://www.cad.zju.edu.cn/home/dengcai/Data/DSH.html>

TABLE 1  
Training and testing time of all algorithms on GIST1M.

Method	Training Time (s)				Test Time (s)			
	$L = 16$	$L = 32$	$L = 64$	$L = 96$	$L = 16$	$L = 32$	$L = 64$	$L = 96$
LSH [8]	0.4	1.0	2.3	2.6	$1.2 \times 10^{-6}$	$2.6 \times 10^{-6}$	$5.8 \times 10^{-6}$	$7.1 \times 10^{-6}$
KLSH [22]	27.4	27.7	27.9	28.3	$30.0 \times 10^{-6}$	$32.3 \times 10^{-6}$	$34.7 \times 10^{-6}$	$36.5 \times 10^{-6}$
SIKH [30]	1.3	2.5	3.8	5.2	$3.9 \times 10^{-6}$	$6.3 \times 10^{-6}$	$10.5 \times 10^{-6}$	$15.9 \times 10^{-6}$
PCAH [34]	31.3	57.2	60.3	75.0	$1.2 \times 10^{-6}$	$2.7 \times 10^{-6}$	$5.6 \times 10^{-6}$	$7.2 \times 10^{-6}$
SpH [35]	42.5	77.8	125.3	239.8	$23.9 \times 10^{-6}$	$42.1 \times 10^{-6}$	$93.4 \times 10^{-6}$	$270.1 \times 10^{-6}$
AGH [24]	340.8	344.7	349.8	356.0	$33.3 \times 10^{-6}$	$52.6 \times 10^{-6}$	$71.2 \times 10^{-6}$	$191.3 \times 10^{-6}$
DSH	33.1	45.9	56.5	63.6	$1.3 \times 10^{-6}$	$2.7 \times 10^{-6}$	$5.8 \times 10^{-6}$	$7.1 \times 10^{-6}$

TABLE 2  
Training and testing time of all algorithms on Flickr1M.

Method	Training Time (s)				Test Time (s)			
	$L = 16$	$L = 32$	$L = 64$	$L = 96$	$L = 16$	$L = 32$	$L = 64$	$L = 96$
LSH [8]	0.3	0.8	1.5	1.8	$0.9 \times 10^{-6}$	$2.0 \times 10^{-6}$	$2.8 \times 10^{-6}$	$4.6 \times 10^{-6}$
KLSH [22]	18.2	18.5	18.9	19.4	$15.2 \times 10^{-6}$	$18.9 \times 10^{-6}$	$22.8 \times 10^{-6}$	$25.1 \times 10^{-6}$
SIKH [30]	1.1	2.0	2.8	4.3	$2.8 \times 10^{-6}$	$3.8 \times 10^{-6}$	$9.1 \times 10^{-6}$	$12.3 \times 10^{-6}$
PCAH [34]	16.7	29.4	31.4	33.7	$1.1 \times 10^{-6}$	$2.1 \times 10^{-6}$	$2.9 \times 10^{-6}$	$4.9 \times 10^{-6}$
SpH [35]	22.3	45.6	106.2	205.5	$16.7 \times 10^{-6}$	$38.5 \times 10^{-6}$	$88.6 \times 10^{-6}$	$251.6 \times 10^{-6}$
AGH [24]	232.9	247.9	257.4	268.1	$28.2 \times 10^{-6}$	$42.2 \times 10^{-6}$	$52.2 \times 10^{-6}$	$155.3 \times 10^{-6}$
DSH	17.4	29.3	35.8	45.9	$0.9 \times 10^{-6}$	$2.1 \times 10^{-6}$	$2.8 \times 10^{-6}$	$4.6 \times 10^{-6}$

TABLE 3  
Training and testing time of all algorithms on SIFT1M.

Method	Training Time (s)				Test Time (s)			
	$L = 16$	$L = 32$	$L = 64$	$L = 96$	$L = 16$	$L = 32$	$L = 64$	$L = 96$
LSH [8]	0.1	0.3	0.6	0.8	$0.4 \times 10^{-6}$	$1.1 \times 10^{-6}$	$1.8 \times 10^{-6}$	$2.4 \times 10^{-6}$
KLSH [22]	10.2	10.4	10.8	11.2	$12.2 \times 10^{-6}$	$13.1 \times 10^{-6}$	$13.8 \times 10^{-6}$	$15.7 \times 10^{-6}$
SIKH [30]	0.5	1.1	2.3	3.5	$0.9 \times 10^{-6}$	$2.3 \times 10^{-6}$	$6.3 \times 10^{-6}$	$7.0 \times 10^{-6}$
PCAH [34]	3.9	6.5	7.5	7.8	$0.5 \times 10^{-6}$	$1.3 \times 10^{-6}$	$2.0 \times 10^{-6}$	$2.5 \times 10^{-6}$
SpH [35]	11.4	28.1	92.7	189.1	$11.8 \times 10^{-6}$	$33.3 \times 10^{-6}$	$77.1 \times 10^{-6}$	$230.9 \times 10^{-6}$
AGH [24]	135.2	142.5	148.1	155.1	$15.3 \times 10^{-6}$	$23.9 \times 10^{-6}$	$31.2 \times 10^{-6}$	$57.1 \times 10^{-6}$
DSH	8.4	12.2	15.5	20.1	$0.5 \times 10^{-6}$	$1.2 \times 10^{-6}$	$1.9 \times 10^{-6}$	$2.6 \times 10^{-6}$

increases. This is consistent with previous work [13], [33] and is probably because that most of the data variance is contained in the top few principal directions so that the later bits are calculated using the low-variance projections, leading to the poorly discriminative codes [33]. By utilizing the geometric structure of the data to guide the projections selection, our DSH successfully combines the advantages of both random projection based methods and the learning based methods. As a result, DSH achieves a satisfied performance on the three data sets and almost outperforms its competitors for all code lengths. It is interesting to see that the performance improvements of DSH over other methods on GIST1M and Flickr1M are larger than that on SIFT1M. Since the dimensions of the data in GIST1M (960-d) and Flickr1M (512-d) are much larger than that in SIFT1M (128-d), this suggests that our DSH method are particularly suitable for high dimensional situations. Figure 3 presents the precision-recall curves of all the algorithms on three data sets with the codes of 48 bits and 96 bits.

Table 1, 2 and 3 show both the training and testing time for different algorithms on three data sets, respectively. We can clearly see that both the training and

testing time of all the methods decrease as the dimension of the data decreases. Considering the training time, the three random projection based algorithms are relatively efficient, especially for LSH and SIKH. KLSH needs to compute a sampled kernel matrix which slows down its computation. The three learning based algorithms are relatively slow, for exploring the data structure. Our DSH is also fast. Although it is slower than the three random projection based algorithms, it is significantly faster than SpH and AGH. Considering the testing time, LSH, PCAH and our DSH are the most efficient methods. All of them simply need a matrix multiplication and a thresholding to obtain the binary codes. SpH consumes much longer time than other methods as the code length increases since it needs to compute the analytical eigenfunctions involving the calculation of trigonometric functions.

### 4.3 Parameter Selection

Our DSH has three parameters:  $p$  (the number of iterations in  $k$ -means),  $\alpha$  (the parameter controlling the groups number) and  $r$  (the parameter for  $r$ -adjacent

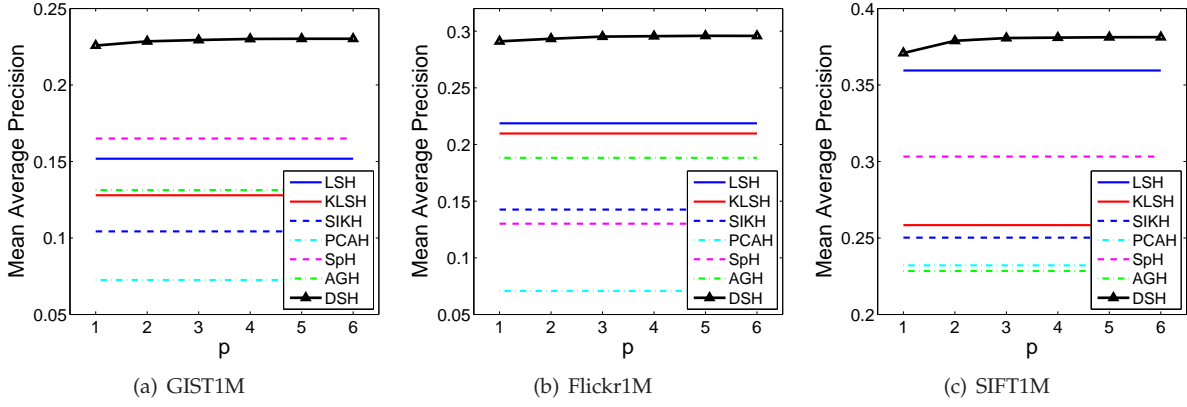


Fig. 4. The performance of DSH vs. the number of iterations of  $k$ -means ( $p$ ) at 64 bits.

TABLE 4  
Training time (s) of DSH vs. the number of iterations of  $k$ -means ( $p$ ) at 64 bits.

Data Set	$p = 1$	$p = 2$	$p = 3$	$p = 4$	$p = 5$	$p = 6$
GIST1M	18.8	37.2	56.5	76.2	94.1	111.7
Flickr1M	11.7	23.5	35.8	48.1	62.6	76.4
SIFT1M	4.8	9.1	15.5	21.2	25.5	31.9

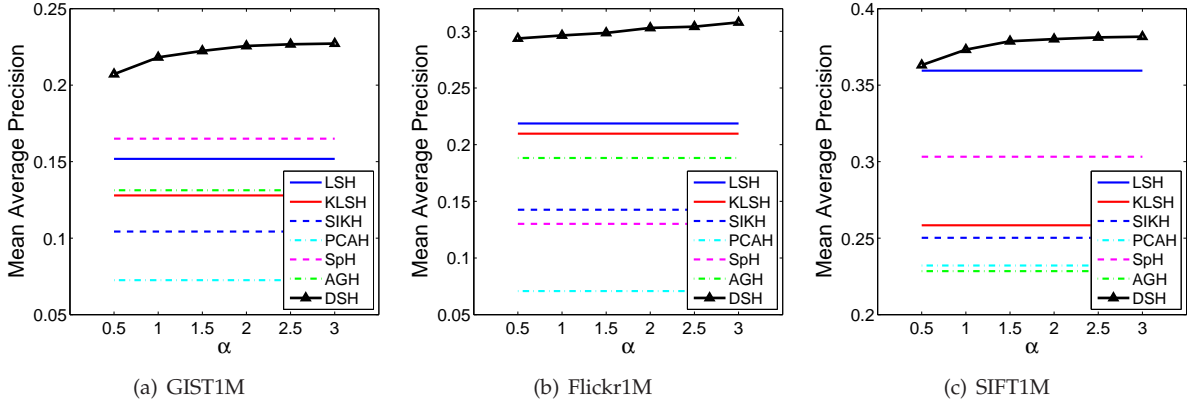


Fig. 5. The performance of DSH vs. the parameter  $\alpha$  (controlling the number of groups) at 64 bits.

TABLE 5  
Training time (s) of DSH vs. the parameter  $\alpha$  (controlling the number of groups) at 64 bits.

Data Set	$\alpha = 0.5$	$\alpha = 1.0$	$\alpha = 1.5$	$\alpha = 2.0$	$\alpha = 2.5$	$\alpha = 3.0$
GIST1M	48.4	52.9	56.5	68.4	74.6	83.5
Flickr1M	21.8	25.3	35.8	46.7	57.3	68.2
SIFT1M	9.8	11.2	15.5	21.3	28.2	37.9

groups). In this subsection, we discuss how the performance of DSH will be influenced by these three parameters. We learn 64-bits hashing codes and the default setting for these parameters is  $p = 3$ ,  $\alpha = 1.5$  and  $r = 3$ . When we study the impact of one parameter, the other parameters are fixed as the default.

Figure 4 and Table 4 show how the performance of DSH varies as the number of iterations in  $k$ -means varies. As the number of iterations increases, it is reasonable to see that both the MAP and the learning time

of DSH increase. On all the three data sets, 3 iterations in  $k$ -means are enough for achieving reasonably good MAP.

Figure 5 and Table 5 show how the performance of DSH varies as  $\alpha$  changes (the groups number generated by  $k$ -means changes). As we can see, as  $\alpha$  becomes larger (the groups number increases), both the MAP and learning time of DSH increase. Setting  $\alpha = 1.5$  is a reasonable balance considering both the accuracy and the efficiency.



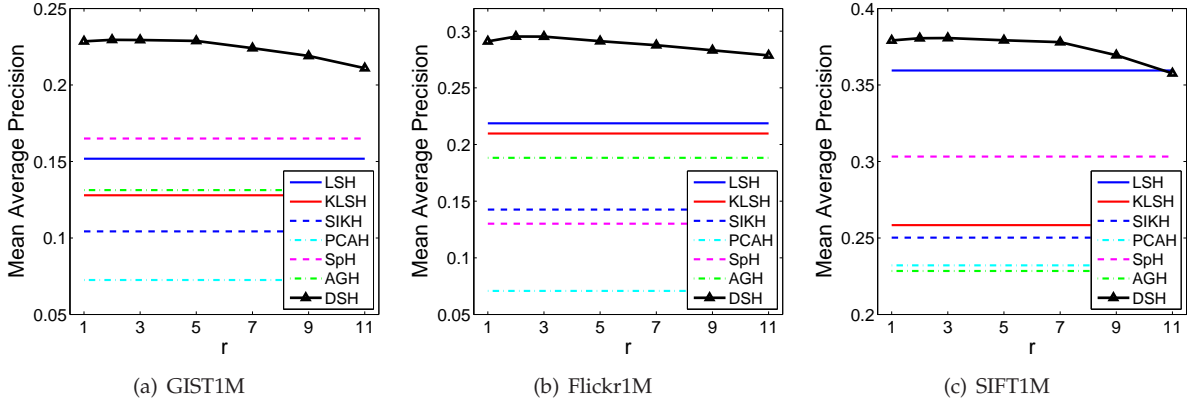


Fig. 6. The performance of DSH vs. the parameter  $r$  (for  $r$ -adjacent groups) at 64 bits.

Figure 6 shows the performance of DSH varies as  $r$  ( $r$ -adjacent groups) changes. DSH achieves stable and consistent good performance as  $r$  is less than 5. As  $r$  becomes larger, DSH generates more projections which are used to separate two far away groups. These projections are usually less critical and redundant. Thus, the performance of DSH decreases.

## 5 CONCLUSION

In this paper, we have developed a novel hashing algorithm, called *Density Sensitive Hashing* (DSH), for high dimensional nearest neighbors search. Different from those random projection based hashing approaches, *e.g.*, Locality Sensitive Hashing, DSH uses the geometric structure of the data to guide the projections selection. As a result, DSH can generate hashing codes with more discriminating power. Empirical studies on three large data sets show that the proposed algorithm scales well to data size and significantly outperforms the state-of-the-art hashing methods in terms of retrieval accuracy.

## REFERENCES

- [1] A. Andoni and P. Indyk. Near-optimal hashing algorithms for approximate nearest neighbor search in high dimensions. *Commun. ACM*, 51(1):117–122, 2008. 1, 2
- [2] L. Arge, M. Berg, H. Haverkort, and K. Yi. The Priority R-tree: a practically efficient and worst-case optimal R-tree. In *SIGMOD*, 2004. 1
- [3] K. B. and T. Darrell. Learning to hash with binary reconstructive embeddings. In *The Neural Information Processing Systems*, 2010. 2
- [4] J. Bentley. Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 18:509–517, 1975. 1
- [5] K. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft. When is nearest neighbor meaningful? In *ICDT*, 1999. 1
- [6] J. Brandt. Transform coding for fast approximate nearest neighbor search in high dimensions. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1815–1822, 2010. 1, 2
- [7] D. Cai. *Spectral Regression: A Regression Framework for Efficient Regularized Subspace Learning*. PhD thesis, Department of Computer Science, University of Illinois at Urbana-Champaign, May 2009. 2
- [8] M. Charikar. Similarity estimation techniques from rounding algorithms. In *STOC*, pages 380–388, 2002. 1, 2, 3, 5, 7
- [9] A. Dasgupta, R. Kumar, and T. Sarls. Fast locality-sensitive hashing. In *IEEE International Conference on Knowledge Discovery and Data Mining*, pages 1073–1081, 2011. 1, 2
- [10] M. Datar, N. Immorlica, P. Indyk, and V. S. Mirrokni. Locality-sensitive hashing scheme based on p-stable distributions. In *Symposium on Computational Geometry 2004*, pages 253–262, 2004. 1, 2
- [11] Y. Gao, B. Zheng, G. Chen, Q. Li, and X. Guo. Continuous visible nearest neighbor query processing in spatial databases. *VLDB J.*, 20(3):371–396, 2011. 1
- [12] A. Gionis, P. Indyk, and R. Motwani. Similarity search in high dimensions via hashing. In *International Conference on Very Large Data Bases*, 1999. 1, 2
- [13] Y. Gong and S. Lazebnik. Iterative quantization: A procrustean approach to learning binary codes. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 817–824, 2011. 2, 6
- [14] J. He, W. Liu, and S.-F. Chang. Scalable similarity search with optimized kernel hashing. In *IEEE International Conference on Knowledge Discovery and Data Mining*, pages 1129–1138, 2010. 1, 2
- [15] H. Jégou, M. Douze, and C. Schmid. Product quantization for nearest neighbor search. *IEEE Trans. Pattern Anal. Mach. Intell.*, 33(1):117–128, 2011. 1, 2
- [16] Q. Jiang and M. Sun. Semi-supervised simhash for efficient document similarity search. In *The 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies*, pages 93–101, 2011. 1, 2
- [17] W. Johnson and J. Lindenstrauss. Extensions of Lipschitz mappings into a Hilbert space. *Contemporary mathematics*, 26:189–206, 1984. 1, 2
- [18] A. Joly and O. Buisson. A posteriori multi-probe locality sensitive hashing. In *ACM Multimedia*, pages 209–218, 2008. 2
- [19] A. Joly and O. Buisson. Random maximum margin hashing. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 873–880, 2011. 1, 2
- [20] M. R. Kolahdouzan and C. Shahabi. Voronoi-based k nearest neighbor search for spatial network databases. In *International Conference on Very Large Data Bases*, pages 840–851, 2004. 2
- [21] F. Korn, N. Sidiropoulos, C. Faloutsos, E. Siegel, and Z. Protopapas. Fast nearest neighbor search in medical image databases. In *International Conference on Very Large Data Bases*, pages 215–226, 1996. 2
- [22] B. Kulis and K. Grauman. Kernelized locality-sensitive hashing for scalable image search. In *IEEE International Conference on Computer Vision*, 2009. 1, 2, 5, 7
- [23] P. Li, A. Shrivastava, J. Moore, and C. Konig. Hashing algorithms for large-scale learning. In *The Neural Information Processing Systems*, 2011. 2
- [24] W. Liu, J. Wang, S. Kumar, and S.-F. Chang. Hashing with graphs. In *International Conference on Machine Learning*, 2011. 2, 6, 7
- [25] Q. Lv, W. Josephson, Z. Wang, M. Charikar, and K. Li. Multi-probe lsh: Efficient indexing for high-dimensional similarity search. In *International Conference on Very Large Data Bases*, pages 950–961, 2007. 2
- [26] Y. Mu, J. Shen, and S. Yan. Weakly-supervised hashing in kernel space. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3344–3351, 2010. 2
- [27] M. Norouzi and D. J. Fleet. Minimal loss hashing for compact binary codes. In *International Conference on Machine Learning*, 2011. 1, 2

- [28] R. Panigrahy. Entropy based nearest neighbor search in high dimensions. In *SODA*, pages 1186–1195, 2006. 2
- [29] L. Paulevé, H. Jégou, and L. Amsaleg. Locality sensitive hashing: A comparison of hash function types and querying mechanisms. *Pattern Recognition Letters*, 31(11):1348–1358, 2010. 2, 3
- [30] M. Raginsky and S. Lazebnik. Locality-sensitive binary codes from shift-invariant kernels. In *The Neural Information Processing Systems*, 2009. 2, 6, 7
- [31] R. Salakhutdinov and G. E. Hinton. Semantic hashing. *International Journal of Approximate Reasoning*, 50(7):969–978, 2009. 1, 2
- [32] J. Wang, O. Kumar, and S.-F. Chang. Semi-supervised hashing for scalable image retrieval. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 3424–3431, 2010. 1, 2, 4
- [33] J. Wang, S. Kumar, and S.-F. Chang. Sequential projection learning for hashing with compact codes. In *International Conference on Machine Learning*, 2010. 2, 5, 6, 7
- [34] X.-J. Wang, L. Zhang, F. Jing, and W.-Y. Ma. Annosearch: Image auto-annotation by search. In *IEEE Conference on Computer Vision and Pattern Recognition*, pages 1483–1490, 2006. 2, 3, 6, 7
- [35] Y. Weiss, A. Torralba, and R. Fergus. Spectral hashing. In *The Neural Information Processing Systems*, pages 1753–1760, 2008. 1, 2, 6, 7
- [36] H. Xu, J. Wang, Z. Li, G. Zeng, N. Yu, and S. Li. Complementary hashing for approximate nearest neighbor search. In *IEEE International Conference on Computer Vision*, 2011. 2, 5
- [37] D. Zhang, J. Wang, D. Cai, and J. Lu. Laplacian co-hashing of terms and documents. In *ECIR*, pages 577–580, 2010. 2
- [38] D. Zhang, J. Wang, D. Cai, and J. Lu. Self-taught hashing for fast similarity search. In *SIGIR*, pages 18–25, 2010. 1, 2, 4